

# Structurally Regularized Non-negative Tensor Factorization for Spatio-temporal Pattern Discoveries

Koh Takeuchi<sup>1,3</sup>, Yoshinobu Kawahara<sup>2,4</sup>, and Tomoharu Iwata<sup>1</sup>

<sup>1</sup> NTT Communication Science Laboratories

{takeuchi.koh, iwata.tomoharu}@lab.ntt.co.jp

<sup>2</sup> The Inst. of Sci. and Ind. Res. (ISIR), Osaka University

ykawahara@sanken.osaka-u.ac.jp

<sup>3</sup> Department of Intelligence Science and Technology, Kyoto University

<sup>4</sup> Center for Advanced Intelligence Project, RIKEN

**Abstract.** Understanding spatio-temporal activities in a city is a typical problem of spatio-temporal data analysis. For this analysis, tensor factorization methods have been widely applied for extracting a few essential patterns into latent factors. Non-negative Tensor Factorization (NTF) is popular because of its capability of learning interpretable factors from non-negative data, simple computation procedures, and dealing with missing observation. However, since existing NTF methods are not fully aware of spatial and temporal dependencies, they often fall short of learning latent factors where a large portion of missing observation exist in data. In this paper, we present a novel NTF method for extracting smooth and flat latent factors by leveraging various kinds of spatial and temporal structures. Our method incorporates a unified structured regularizer into NTF that can represent various kinds of auxiliary information, such as an order of timestamps, a daily and weekly periodicity, distances between sensor locations, and areas of locations. For the estimation of the factors for our model, we present a simple and efficient optimization procedure based on the alternating direction method of multipliers. In missing value interpolation experiments of traffic flow data and bike-sharing system data, we demonstrate that our proposed method improved interpolation performances from existing NTF, especially when a large portion of missing values exists.

## 1 Introduction

Spatio-temporal data covering a wide area of a city have become available due to the commoditization of sensor-monitoring systems and mobile-phone networks. These monitoring systems observe various types of data, such as vehicle transportation counts on a road network, bike-renting counts of a bike-sharing system, and the purchasing records of shops around a city, where missing values often appear due to the failure of sensor nodes, data transmission errors, and trouble with data recording systems. We can find rich and bounteous information in such spatio-temporal data. However, it becomes difficult to grasp what spatio-temporal activities appeared in the data at a glance. Therefore, understanding of such activities via pattern extractions is a typical problem of spatio-temporal data analysis, in which the interpretability of the extracted patterns is regarded as one of the most important property for analysis methods.

Tensor factorization methods have been widely applied to discover spatial and temporal patterns from various kinds of spatio-temporal data [17]. These methods represent spatio-temporal data as a higher-order dimensional array, called a tensor that is a generalization of a matrix. For example, we can represent spatio-temporal data as a three-way tensor whose first, second, and third modes correspond to sensor locations, timestamps for 24 hours, and the observed days. We illustrated an example of a tensor for spatio-temporal data analysis in Fig. 1. With this formulation, we can naturally incorporate an assumption that daily or weekly periodicity can be found in data and similar spatial patterns appear on different days. We can extract a few numbers of spatial, temporal, and daily patterns into latent factors by decomposing the tensor. However, since most existing tensor factorization methods do not consider the non-negativity of data where observations only contain non-negative values, they often result in messy and hard to interpret factors.

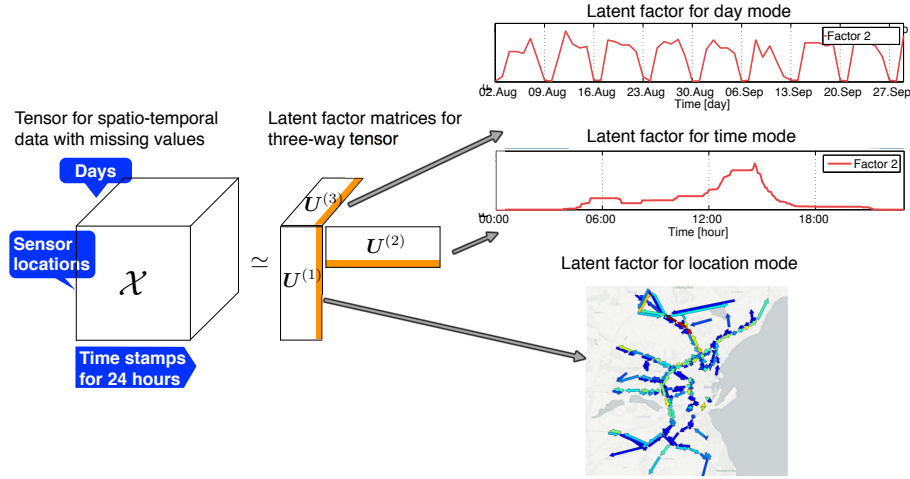


Fig. 1: Example for a non-negative tensor factorization method on analysis of a traffic-flow data set, where latent patterns for location, time, and day modes are extracted in the latent factor.

Unlike those tensor factorization methods, Non-negative Tensor Factorization (NTF) [8], which leverages non-negativity, is effective for extracting interpretable patterns from the non-negative data [18, 13]. This method has successively yielded interpretable factors from various kinds of spatio-temporal data, such as location-based social network services [25, 14], mobile phone GPS logs [10], log messages of network equipment [16], and traffic records of road networks [32]. However, NTF was not applicable to the existence of missing values. To deal with missing values, NTF was recently extended to learn the latent factors from a subset of elements in a tensor, called the non-negative tensor completion [31, 15]. With this NTF, we can interpolate missing values in data by learned latent factors. However, NTF methods for the missing value completion problem suffer from overfitting when just a few observations are available. Because they

ignore spatial and temporal contextual information such as the order of time stamps, weekly periodicity, the distances between sensor locations and treats each feature of the tensor independently.

To incorporate such contextual information, most matrix/tensor factorization methods have employed a graph Laplacian based regularizer for encouraging the latent factors to be smooth with spatio-temporal dependencies [21]. The graph regularized non-negative matrix factorization [6] is a variant of such schemes and has been widely utilized in many applications, however, it does not consider scenarios where missing values exist and analyzing higher-order dimensional arrays.

Another choice for representing such auxiliary information is structured regularizers [2] that have become popular in the fields of machine learning, signal processing, and data mining [7, 28]. For example, the fused lasso [27], which is also known as the total variation, approximates parameters by piecewise-constant values with the order of parameters. Since its estimated parameters have the same estimated value, this is beneficial for finding segments of parameters. In a pioneering work [29], the penalized matrix decomposition was proposed to utilize the fused lasso as a regularizer on latent factors and was applied to a gene data analysis problem. They presented latent factors easy to find gene segments rather than existing matrix factorization methods incorporated the lasso regularizer. However, this method and its subsequent works have only considered the fused lasso without incorporating more general structured regularizer such as spatial dependencies of sensors, and also ignored the non-negative properties and the existence of missing values.

In this paper, we attempt to solve a problem of extracting latent factors from spatio-temporal data where a lot of missing values exists. To tackle this problem, we propose a novel NTF that learns factors by employing spatial and temporal auxiliary information as regularizers. We utilize this information to represent phenomena often appear in spatio-temporal data, such as counts of vehicles passed roads smoothly grow or decrease or take the same value along with space and time. To exploit such information, we introduce a regularizer that consists of both a graph-based Laplacian regularizer and structured regularizers that incorporate not only the order of features but also more general graph and group based structures [3, 24]. With our regularizer, we can utilize various kinds of auxiliary information into NTF including a daily and weekly periodicity, distances between sensor locations, and areas of locations. Our proposed method is highly robust to the presence of a large portion of missing values because it encourages latent factors to be smooth and flat with spatial and temporal structures, where we regard segments of parameters that take the same value as flat. To estimate the latent factors for our proposed method, we present an efficient optimization procedure of the alternating direction method of multipliers [4] that utilizes simple proximity operators of the conjugate gradient method [21] and a parametric network flow algorithm [12].

We conducted missing value interpolation experiments with real-world traffic flow data and compared the performance of our proposed method with existing NTF methods. We demonstrate that our proposed method improved the interpolation performances from existing NTF methods. We also show that our extracted factors were interpretable to detect change points. Because our factors have segments, we can easily find a boundary of segments as a change point.

## 2 Non-negative Tensor Factorization

We denote a  $N$ -th way non-negative tensor as  $\mathcal{X} \in \mathbb{R}_{\geq 0}^{I_1 \times \dots \times I_N}$ , where  $I_n$  is the number of features in the  $n$ -th mode. The  $n$ -th mode unfolding of a tensor  $\mathcal{X}$  is denoted as  $\mathcal{X}_n$ . We use  $i = (i_1, \dots, i_N)$  and  $D$  to represent an element and the whole set of the elements in the tensor, respectively. A subset of the observed elements in the tensor is denoted by  $\Omega = \{i \mid x_i \text{ is observed}, \forall i \in D\}$ .

NTF decomposes the observed values of tensor  $\mathcal{X}$  into  $K$  latent non-negative factors, where  $K \ll \min(I_1, \dots, I_N)$ . The  $n$ -th mode factor matrix is denoted as  $\mathbf{A}^{(n)} \in \mathbb{R}_{\geq 0}^{I_n \times K}$  whose  $k$ -th column is factor vector  $\mathbf{a}_k^{(n)} \in \mathbb{R}_{\geq 0}^{I_n}$ . We denote a whole set of factor vectors as  $\mathbf{A} = \{\mathbf{a}_k^{(n)} \mid \forall (n, k)\}$ . An estimation for element  $x_i$  is given by a sum of latent factor vectors  $\hat{x}_i = \sum_{k=1}^K a_{i_1,k}^{(1)} a_{i_2,k}^{(2)} \dots a_{i_N,k}^{(N)} \in \hat{\mathcal{X}}$ . We denote the transpose operator as  $\top$ , the Khatri-Rao product as  $\odot$ , and its series as  $\odot_{n=1}^N \mathbf{A}^{(n)} = \mathbf{A}^{(1)} \odot \dots \odot \mathbf{A}^{(N)}$ .

The empirical loss function for NTF can be defined as a sum of divergences that indicates a discrepancy between  $x_i$  and its estimation  $\hat{x}_i$ :

$$f(A) = D_{\Omega}(\mathcal{X} \parallel \hat{\mathcal{X}}) + \sum_{n=1}^N \sum_{k=1}^K g^{(n)}(\mathbf{a}_k^{(n)}), \quad (1)$$

where  $D_{\Omega}(\mathcal{X} \parallel \hat{\mathcal{X}}) = \sum_{i \in \Omega} d(x_i \parallel \hat{x}_i)$ . We  $d(p \parallel q)$  to denote a divergence between scalars  $p$  and  $q$ , and  $g^{(n)}$  to denote a penalty function for the  $n$ -th mode factor vector. Because loss function  $f$  is non-convex with respect to  $A$ , an NTF problem is to obtain a local minimizer  $A^*$  of the loss under a non-negative constraint:

$$A^* = \arg \min_A f(A) \text{ subject to } \mathbf{a}_k^{(n)} \geq 0 \quad \forall (n, k). \quad (2)$$

The graph regularized non-negative matrix factorization method [6] employs a graph Laplacian regularizer [22] to represent the smoothness in latent factors. An adjacency matrix for the  $n$ -th mode features is denoted as  $\mathbf{W}^{(n)} \in \mathbb{R}^{I_n \times I_n}$  that represents a graph whose nodes and capacities of edges correspond to the features of the  $n$ -th mode and the similarity measures between the two features, respectively. The Laplacian matrix can be denoted as  $\mathbf{L}^{(n)} = \mathbf{D}^{(n)} - \mathbf{W}^{(n)}$ , where  $\mathbf{D}^{(n)}$  is a diagonal matrix whose elements are the sums of each row of  $\mathbf{W}^{(n)}$ . Then a graph Laplacian regularizer can be defined:

$$g^{(n)}(\mathbf{a}_k^{(n)}) = \mathbf{a}_k^{(n)\top} \mathbf{L}^{(n)} \mathbf{a}_k^{(n)}. \quad (3)$$

This regularizer penalty function encourages smoothness because its formulation equals putting a weighted quadratic term on the difference between the adjacency elements.

## 3 Proposed model

We introduce a unified structured regularizer to employ both smooth and piecewise-constant properties with auxiliary structures:

$$g^{(n)}(\mathbf{a}_k^{(n)}) = \sum_{m=1}^3 \lambda_m g_m^{(n)}(\mathbf{a}_k^{(n)}) + g_{\geq 0}^{(n)}(\mathbf{a}_k^{(n)}), \quad (4)$$

where  $\lambda_1, \lambda_2$  and  $\lambda_3$  are the hyperparameters for each regularizer. We employ a Generalized Fused Lasso (GFL) [5, 30] and a Higher-Order Fused Lasso (HOFL) [24] as  $g_1^{(n)}$  and  $g_2^{(n)}$ , respectively.  $g_3^{(n)}$  corresponds to the Laplacian regularizer for extracting smooth patterns. We use an indicator function for the non-negative region:

$$g_{\geq 0}^{(n)}(\mathbf{a}_k^{(n)}) = \begin{cases} 0 & (\text{if } a_{i,k} \geq 0, \forall i) \\ +\infty & (\text{otherwise}) \end{cases}. \quad (5)$$

The GFL penalty is defined:

$$g_1(\mathbf{a}_k^{(n)}) = \sum_{j=1}^{I_n} \sum_{j'=1}^{I_n} w_{j,j'}^{(n)} |a_{j,k}^{(n)} - a_{j',k}^{(n)}|. \quad (6)$$

The GFL prefers parameters with the same value if they are adjacent on the given graph, such as distances between sensor locations and temporal lags between time stamps. The HOFL encourages parameters in a given group to take identical values [24]. With this regularizer, we can utilize auxiliary information, such as sensors placed in a specific area that may output similar values and a group of time stamps when a specific train leaves from a station. We denote the  $r$ -th group of features in the  $n$ -th mode as  $g_r^{(n)} \subseteq D_n$  and a set of groups by  $\mathcal{G}^{(n)} = \{g_1^{(n)}, \dots, g_{R_n}^{(n)}\}$ , where  $D_n$  and  $R_n$  are a set of elements in the  $n$ -th mode and the number of groups, respectively. The weights of each element for the  $r$ -th group on the  $n$ -th mode are denoted by  $c_{r,m}^{(n)} = \bar{c}_{r,m}^{(n)}$  if  $m \in g_r^{(n)}$ , and 0 otherwise, where  $\bar{c}_{r,m}^{(n)} > 0$ . Then a simplified HOFL penalty  $g_2(\mathbf{a}_k^{(n)})$  is given:

$$\sum_{r=1}^R \sum_{m=1}^{I_n} c_{r,j_m}^{(n)} |a_{j_m,k}^{(n)} - \bar{a}_{r,j_m,k}^{(n)}| + \theta_r^{(n)} (a_{s_r,k}^{(n)} - a_{t_r,k}^{(n)}), \quad (7)$$

where  $\theta_r^{(n)} > 0$  is a hyperparameter that controls the consistency of the parameters in a group.  $\bar{a}_{r,k}^{(n)}$  is defined as  $\bar{a}_{r,m,k}^{(n)} = a_{s_k,k}^{(n)}$  (if  $m \geq s_k$ ),  $a_{t_k,k}^{(n)}$  (if  $m \leq t_k$ ) and  $a_{j_m,k}^{(n)}$  (otherwise) for distinct indices  $j_1, j_2, \dots, j_{I_n} \in D_n$  that correspond to a permutation that arranges the entries of  $\mathbf{a}_k^{(n)}$  in a non-increasing order. Thresholding indices  $s_r$  and  $t_r$  are given as  $s_k = \min \{m' \mid \sum_{m=1}^{m'} c_{r,j_m}^{(n)} \geq \theta_r^{(n)}\}$  and  $t_k = \min \{m' \mid \sum_{m=m'}^{I_n} c_{r,j_m}^{(n)} < \theta_r^{(n)}\}$ .

For convenience, we denote  $\bar{g}^{(n)}(\mathbf{a}_k^{(n)}) = \sum_{m=1}^2 \lambda_m g_m^{(n)}(\mathbf{a}_k^{(n)}) + g_{\geq 0}^{(n)}(\mathbf{a}_k^{(n)})$ . By adopting our structured regularizers to the loss of NTF, we define the following minimization problem for our purpose:

$$A^* = \arg \min_A D_\Omega(\mathcal{X} \|\hat{\mathcal{X}}) + \sum_{n=1}^N \sum_{k=1}^K \bar{g}^{(n)}(\mathbf{a}_k^{(n)}) + \lambda_3 g_3^{(n)}(\mathbf{a}_k^{(n)}). \quad (8)$$

Note that when  $\lambda_1 = \lambda_2 = \lambda_3 = 0$ , our method is reduced to an original NTF. When  $\lambda_2 = \lambda_3 = 0$ , our method can be regarded as a tensor extension of the graph-regularized non-negative matrix factorization. Our method includes those methods as special cases.

## 4 Parameter estimation

We present an efficient parameter estimation procedure for obtaining a local minimizer of our proposed method. We employ a scaled formulation of the Alternating Direction Method of Multipliers (ADMM) for NTF [15]. The minimization problem for our proposed method can be rewritten:

$$\begin{aligned} \min_{A, \mathcal{Z}} D_{\Omega}(\mathcal{X} \| \mathcal{Z}) + \sum_{n=1}^N \sum_{k=1}^K \bar{g}^{(n)}(\mathbf{b}_k^{(n)}) + g_3^{(n)}(\mathbf{a}_k^{(n)}) \\ \text{subject to } \mathcal{Z} = \hat{\mathcal{X}}, \mathbf{a}_k^{(n)} = \mathbf{b}_k^{(n)} \ (\forall n, k), \end{aligned} \quad (9)$$

where  $\mathcal{Z}$  and  $\mathbf{b}_k^{(n)}$  are auxiliary variables, and  $P_{\Omega}$  is a projection function that only retains the divergence of the observed elements. To solve our problem efficiently with keeping both constraints and separability, we define an augmented Lagrangian for our problem:

$$\begin{aligned} L_{\rho}(A, B, \mathcal{Z}) = D_{\Omega}(\mathcal{X} \| \mathcal{Z}) + \frac{\rho}{2} \|\mathcal{Z} - \hat{\mathcal{X}} + \mathcal{U}\|_{\mathcal{F}}^2 \\ \sum_{n=1}^N \sum_{k=1}^K \bar{g}^{(n)}(\mathbf{b}_k^{(n)}) + g_3^{(n)}(\mathbf{a}_k^{(n)}) + \frac{\rho}{2} \|\mathbf{a}_k^{(n)} - \mathbf{b}_k^{(n)} + \mathbf{u}_k^{(n)}\|_2^2, \end{aligned} \quad (10)$$

where  $\mathcal{U}$  and  $\mathbf{u}_k^{(n)}$  are Lagrangian multipliers, and  $\rho$  is a step-size parameter, respectively. We summarize the minimization procedure for our proposed method in Algorithm 1. The minimization for ADMM can be efficiently calculated if a simple minimization operator for each of each  $\mathbf{a}_k^{(n)}$  and  $\mathbf{b}_k^{(n)}$  exists.

The loss function with respect to  $\mathbf{A}^{(n)}$  and  $\mathbf{b}_k^{(n)}$  contains the graph Laplacian regularizer and the non-separable graph-based and group-based penalties, respectively. Thus the main difficulty with our proposed method lies in the minimization of  $\mathbf{A}^{(n)}$  and  $\mathbf{b}_k^{(n)}$ , whose minimization problems can be rewritten:

$$\mathbf{A}^{(n)} = \arg \min_{\mathbf{A}^{(n)}} \frac{\rho}{2} \|\bar{\mathcal{Z}}_n - \mathbf{A}^{(n)} \mathbf{V}_n^{\top}\|_2^2 + \frac{\rho}{2} \|\mathbf{A}^{(n)} - \bar{\mathbf{V}}_n\|_2^2 + \lambda_3 \sum_{k=1}^K g_3^{(n)}(\mathbf{a}_k^{(n)}) \quad (11)$$

$$\mathbf{b}_k^{(n)} = \arg \min_{\mathbf{b}_k^{(n)}} \bar{g}^{(n)}(\mathbf{b}_k^{(n)}) + \frac{\rho}{2} \|\bar{\mathbf{v}}_k^{(n)} - \mathbf{b}_k^{(n)}\|_2^2, \quad (12)$$

where  $\bar{\mathcal{Z}} = \mathcal{Z} + \mathcal{U}$ ,  $\mathbf{V}_n = \odot_{n=n'}^N \mathbf{A}^{(n)}$ ,  $\bar{\mathbf{V}}_n = \mathbf{B}^{(n)} - \mathbf{U}^{(n)}$ , and  $\bar{\mathbf{v}}_k^{(n)} = \mathbf{a}_k^{(n)} + \mathbf{u}_k^{(n)}$ . We efficiently solve the minimization of Eq. (11) by using the fact that it corresponds to the loss function of the graph regularized alternating least squares [21], which approximately runs in  $\mathcal{O}(\text{nnz}(\mathbf{L}^{(n)})K)$  ( $\text{nnz}(\cdot)$  is the number of non-zero elements).

The minimization problem in Eq. (12) corresponds to the calculation of the proximity operator, which is defined as:  $\text{prox}_{\gamma h}(\theta) = \arg \min_{\theta} h(\theta) + \frac{1}{2\gamma} \|\hat{\theta} - \theta\|_2^2$ . We present a minimization procedure for Eq. (12) by leveraging the properties of the proximity operator, and obtaining a minimizer for the sum of the non-negative indication function and other convex functions by the following property [26]:  $\text{prox}_{\bar{g}^{(n)}} =$

---

**Algorithm 1:** Alternative direction method of multiplier for our proposed non-negative tensor factorization

---

**Input :**  $\mathcal{X}, \Omega, \lambda_1, \lambda_2, \lambda_3, K, \mathbf{W}^{(n)}$   
**Output:** set of factor matrices  $A$

- 1 Initialize parameters
- 2 Sample  $A, B$ , and  $\mathcal{Z}$  from random distributions
- 3 **repeat**
- 4     Alternatively update parameters;
- 5      $\mathcal{Z} \leftarrow \arg \min_{\mathcal{Z}} D_{\Omega}(\mathcal{X} \parallel \mathcal{Z}) + (\rho/2) \|\mathcal{Z} - \hat{\mathcal{X}} + \mathcal{U}\|_F^2$
- 6     **for**  $n = 1$  **to**  $N$  **do**
- 7         Update  $A^{(n)}$  by solving Eq. (11)
- 8         **for**  $k = 1$  **to**  $K$  **do**
- 9             Update  $b_k^{(n)}$  by solving Eq. (12)
- 10         **end**
- 11          $U^{(n)} \leftarrow U^{(n)} + (A^{(n)} - B^{(n)})$
- 12     **end**
- 13      $U \leftarrow U + (\mathcal{X} - \hat{\mathcal{X}})$
- 14 **until** *convergence*;

---

$\text{prox}_{g_{\geq 0}^{(n)}} \circ \text{prox}_{\lambda_1 g_1^{(n)} + \lambda_2 g_2^{(n)}}$ . Thus, if we have a minimizer for  $\lambda_1 g_1^{(n)} + \lambda_2 g_2^{(n)}$ , we can attain the exact minimizer for  $\bar{g}^{(n)}$  by setting negative parameters to zeros. A minimizer for  $\lambda_1 g_1^{(n)} + \lambda_2 g_2^{(n)}$  can be simply calculated by employing a submodular function minimization procedure. Because the penalty functions of GFL and HOFL are the Lovász extensions [19] of the graph-representable submodular functions [11], we can attain a minimizer for the sum of functions  $\lambda_1 g_1^{(n)} + \lambda_2 g_2^{(n)}$  by an efficient parametric network flow algorithm [7, 30, 24]. We show the details of our minimization procedure for this function in the appendix.

## 5 Related Works

There has been a lot of articles in which NTF was applied to analyze spatio-temporal data. Kimura et. al proposed a special NTF that decomposes a three-way tensor into two-factor matrices and a three-mode tensor for extracting log messages related to network failures [16]. Yang et. al proposed a combination of NTF without regularizers and post-processing for modeling user activities [32]. Koh et. al proposed an NTF that simultaneously decomposes multiple tensors to extract patterns appeared among different tensors[25]. NTF was used to extract spatio-temporal patterns from human-flow data [10]. However, all of those methods did not employ regularizers into NTF and their methods were not applicable to missing values. One exception is a paper of Sun et. al [23], in which they proposed a probabilistic non-negative Tucker decomposition for discovering interactions among factors. However, they did not incorporate the spatial and temporal structures into regularizers. Han et. al proposed an extension of NTF for predicting future observations [14]. However, they did not consider spatial struc-

tures. Our method can be applied to their framework to utilizing spatial and temporal regularizers. The estimation procedures of them were based on the multiplicative update rule and EM algorithm. Our proposed method can utilize graphs and groups of spatial and temporal features to regularize parameters and also employ ADMM as an estimation procedure.

## 6 Experiments

We conducted missing value completion problems with a traffic flow data set provided by City Pulse [1] and two bike-sharing system data sets recorded in Washington D. C.<sup>5</sup> and New York<sup>6</sup> [1].

The traffic flow data consist of the numbers of cars that passed at 419 locations every thirty minutes in Arhus City, Denmark. We picked 30 days from August 2nd to 31st 2014, and constructed three-way tensor  $\mathcal{X} \in \mathbb{R}^{48 \times 30 \times 441}$  whose modes corresponded to 48 daily time points, 30 days, and 441 observation locations, respectively. From the bike-sharing system data in Washington D.C. and New York, we employed 15 days from April 1st to the 15th with 351 and 344 bike stations. We constructed three-way tensors  $\mathcal{X} \in \mathbb{R}^{24 \times 15 \times 351}$  and  $\mathcal{X} \in \mathbb{R}^{24 \times 15 \times 344}$  whose values were the numbers of bikes returned to the station in an hour. For the time mode, we utilized the adjacency of the time points as a graph. For the day mode, we employed the adjacency of days and the days of the week as a graph and groups. For the location mode, we used the inverse of the Euclid distance of GPS locations and clusters attained by k-nearest neighbors ( $k = 5, 10$ ) for a graph and groups.

We exploited the Euclid distance as the divergence in experiments. We compared our proposed method (Proposed 1) and our proposed method with only the graph Laplacian regularizer (Proposed 2,  $\lambda_1 = \lambda_2 = 0$ ) with NTF estimated by ADMM [15] (ADMM), NTF with the graph Laplacian regularizer [6] estimated by a multiplicative update rule considering missing values (Multi+Lap) [9], and NTF estimated by the multiplicative update rule (Multi). We set the proportion of observations to  $p = \{0.1, 0.01, 0.005, 0.001\}$ . By five-fold cross validation, we selected  $K$  and other hyperparameters from  $K = \{3, 5, 10\}$  and  $\{0.1, 1, 10\}$ . We utilized the normalized RMSE (NRMSE) and the normalized deviation (ND) as error measurements:

$$\text{NRMSE} = \sqrt{(1/|\Omega|) \sum_{(p,t) \in \Omega} (x_{p,t} - \hat{x}_{p,t})^2 / Q}, \quad (13)$$

$$\text{ND} = (1/|\Omega|) \sum_{(p,t) \in \Omega} |x_{p,t} - \hat{x}_{p,t}| / Q, \quad (14)$$

here  $Q = (1/|\Omega|) \sum_{(p,t) \in \Omega} |x_{p,t}|$ . We ran our experiments five times with randomly selected different missing values.

The results are shown in Tables 1, 2, 3, 4, 5, and 6, where the left and right values in a cell correspond to the average and the standard deviation of those values. We

<sup>5</sup> <https://www.capitalbikeshare.com>

<sup>6</sup> <http://www.citibikenyc.com/>

Table 1: NNAME for the traffic flow data of our proposed method (Proposed 1), our proposed method with the graph Laplacian regularizer (Proposed 2), NTF estimated by ADMM (ADMM), NTF with the graph Laplacian regularizer estimated by the multiplicative update rule (Multi+Lap), and NTF (Multi)

	$p = 0.1$	$p = 0.01$	$p = 0.005$	$p = 0.001$
Proposed 1	<b>0.50 (0.00)</b>	<b>0.99 (0.03)</b>	<b>1.49 (0.03)</b>	<b>1.87 (0.01)</b>
Proposed 2	0.51 (0.00)	1.12 (0.05)	<b>1.49 (0.03)</b>	1.89 (0.01)
ADMM	0.51 (0.00)	1.15 (0.03)	1.50 (0.02)	1.91 (0.00)
Multi+Lap	0.52 (0.00)	2.98 (1.86)	2.92 (1.00)	11.9 (11.7)
Multi	0.52 (0.00)	2.89 (2.22)	3.27 (1.80)	5.98 (6.24)

Table 2: NNAME for the bike-sharing record data of Washington D.C.

Method	$p = 0.1$	$p = 0.01$	$p = 0.005$	$p = 0.001$
Proposed 1	<b>1.67 (0.02)</b>	<b>2.14 (0.02)</b>	<b>2.21 (0.04)</b>	<b>2.43 (0.02)</b>
Proposed 2	1.68 (0.01)	<b>2.14 (0.05)</b>	2.22 (0.05)	2.62 (0.08)
ADMM	1.68 (0.02)	2.21 (0.05)	2.32 (0.03)	2.47 (0.01)
Multi+Lap.	1.69 (0.01)	2.72 (0.22)	2.76 (0.24)	11.2 (4.62)
Multi	1.70 (0.01)	299.1 (405.7)	8.25 (4.87)	16.3 (3.13)

confirmed that our proposed methods showed the best performance in every setting. Our proposed method was robust to the appearance of a large portion of missing values for every data set  $p = \{0.01, 0.005, 0.001\}$ . Our proposed method with both the graph-based Laplacian and structured regularizer (Proposed 1) showed better or competitive performance with our proposed method with the graph-based Laplacian regularizer (Proposed 2). Furthermore, our proposed method with the graph-based Laplacian regularizer (Proposed 2) always outperformed the same model estimated by the multiplicative update rule (Multi+Lap). This result was caused by the benefits of simultaneously combining graph-based and structured regularizers with graph and group structures. Thus our proposed model and parameter estimation procedure both contributed to the improvements on missing value interpolations. The existing methods resulted in poor performances with settings where a large portions of tensor elements were missing.

To check the qualitative performances of the interpretability, we showed the extracted factors of proposed method (Proposed 1) and existing NTF with the Laplacian regularizer (Multi+Lap) from traffic flow data in Figs. 2, 3, 4, 5, 6, and 7, where  $p = 0.1$ . The degree of freedom (DoF) in Figures corresponded the number of segments in a factor matrix. Thanks to the Laplacian and structured regularizers, proposed method extracted the interpretable latent factors in which both smooth and flat properties appeared, whose DoF of parameters in factor matrices were extremely less than that of NTF with the Laplacian regularizer. Our factors with low DoF were easy to find change points. For example, the blue factor had a change at 3 am and gradually grew until 6 am. Then it took the constant values until 3 pm in Fig. 2. This factor also has the same value from day 2 to day 6 and from day 8 to day 13. Thus, we can easily understand that the blue factor in Figs. 2 and 4 corresponded to activity that occurred in weekday during

Table 3: NRAM for the bike-sharing record data of New York

Method	$p = 0.1$	$p = 0.01$	$p = 0.005$	$p = 0.001$
Proposed 1	<b>0.98 (0.00)</b>	<b>1.28 (0.02)</b>	<b>1.42 (0.01)</b>	<b>1.62 (0.01)</b>
Proposed 2	<b>0.98 (0.00)</b>	1.30 (0.03)	1.44 (0.01)	1.63 (0.01)
ADMM	<b>0.98 (0.00)</b>	1.34 (0.02)	1.49 (0.01)	1.65 (0.01)
Multi+Lap	1.00 (0.02)	27.2 (41.3)	5.68 (3.22)	1.86 (0.17)
Multi	1.00 (0.02)	53.7 (22.0)	26.6 (29.4)	3.04 (1.03)

Table 4: ND for the traffic flow data of our proposed method (Proposed 1), our proposed method with the graph Laplacian regularizer (Proposed 2), NTF estimated by ADMM (ADMM), NTF with the graph Laplacian regularizer estimated by the multiplicative update rule (Multi+Lap), and NTF (Multi)

Method	$p = 0.1$	$p = 0.01$	$p = 0.005$	$p = 0.001$
Proposed 1	<b>0.27 (0.00)</b>	<b>0.46 (0.01)</b>	<b>0.70 (0.01)</b>	<b>0.92 (0.01)</b>
Proposed 2	0.28 (0.00)	0.51 (0.02)	0.71 (0.02)	0.94 (0.00)
ADMM	0.28 (0.00)	0.53 (0.01)	0.73 (0.01)	0.94 (0.00)
Multi+Lap	0.28 (0.00)	0.61 (0.08)	0.78 (0.02)	1.19 (0.16)
Multi	0.28 (0.00)	0.60 (0.08)	0.81 (0.04)	1.18 (0.24)

daylight with a spatial pattern in Fig. 6. However, NMF with the Laplacian regularizer resulted in messy factors. We also showed that of bike-sharing data in Washington D.C. in Figs. 8, 9, 10, 11, 12, and 13. Our proposed method also extracted more interpretable patterns than existing NTF. For example, the yellow factor of ours in Fig. 8 had a change point at 8 am. After it had taken a peak at 12 am, it kept the same value from 1 pm to 5 pm. Then its value gradually decreased to zero. The yellow factor in Fig. 10 had the same high value on day 2, 3, 9, and 10. Thus, we confirmed that this factor indicated a weekend afternoon activity with a spatial pattern in Fig. 12. Similar interpretations can be obtained from other factors of ours.

## 7 Conclusion

In this paper, we proposed a structurally regularized non-negative tensor factorization that incorporated both the graph Laplacian and the structured regularizers on latent factors. For the structured regularizer, we employed the generalized fused lasso and the higher-order fused lasso to represent both graph-based and group-based information in time and space. We introduced a flexible and efficient parameter estimation method based on the alternating direction method of multipliers and showed a proximity operator for our unified structured regularizer. With experiments on a missing value imputation problem of three data sets, we confirmed that our proposed method showed the best quantitative performance and successfully extracted more interpretable latent factors than the existing non-negative tensor factorization methods.

**Acknowledgements** The part of this work was supported by JSPS KAKENHI Grant Numbers JP16H01548 and JP26280086, and NICT "Research and Development on Fundamental and Utilization Technologies for Social Big Data".

Table 5: ND for the bike-sharing record data of Washington D.C.

Method	$p = 0.1$	$p = 0.01$	$p = 0.005$	$p = 0.001$
Proposed 1	<b>0.81 (0.00)</b>	<b>0.91 (0.01)</b>	<b>1.04 (0.02)</b>	<b>1.04 (0.08)</b>
Proposed 2	<b>0.81 (0.01)</b>	<b>0.91 (0.00)</b>	<b>1.04 (0.02)</b>	1.05 (0.09)
ADMM	<b>0.81 (0.01)</b>	<b>0.91 (0.01)</b>	1.10 (0.01)	1.20 (0.01)
Multi+Lap	<b>0.81 (0.00)</b>	1.19 (0.05)	1.51 (0.15)	1.99 (0.54)
Multi	<b>0.81 (0.00)</b>	9.70 (6.91)	1.44 (0.10)	2.56 (0.31)

Table 6: ND for the bike-sharing record data of New York

Method	$p = 0.1$	$p = 0.01$	$p = 0.005$	$p = 0.001$
Proposed 1	<b>0.60 (0.00)</b>	<b>0.72 (0.01)</b>	<b>0.81 (0.01)</b>	<b>0.93 (0.01)</b>
Proposed 2	<b>0.60 (0.00)</b>	0.73 (0.01)	0.82 (0.01)	<b>0.93 (0.01)</b>
ADMM	<b>0.60 (0.00)</b>	0.74 (0.01)	0.84 (0.01)	0.94 (0.01)
Multi+Lap	<b>0.60 (0.00)</b>	2.13 (1.48)	1.53 (0.21)	1.05 (0.05)
Multi	<b>0.60 (0.00)</b>	4.48 (0.79)	2.51 (0.95)	1.21 (0.07)

## A Appendix

Although the issue in Eq. (12) is a general problem containing the previous problems [5, 30, 24] as special cases, we can solve it in a similar manner as these works. We first briefly introduce the parametric optimization method for a non-decreasing set function. Let  $\alpha \in \mathbb{R}_{\geq 0}$ , and define set function  $l_\alpha(S) = l(S) - \alpha \mathbf{1}(S)$  ( $\forall S \subset V$ ), where  $\mathbf{1}(S) = \sum_{i \in S} 1$ . Then if  $l$  is a non-decreasing submodular function, then there exists a set of  $r + 1$  ( $\leq |V|$ ) subsets:  $S^* = \{S_0 \subset S_1 \subset \dots \subset S_r\}$ , where  $S_j \subset V$ ,  $S_0 = \emptyset$ , and  $S_r = V$ , and  $r + 1$  subintervals  $Q_r$  of  $\alpha$ :  $Q_0 = [0, \alpha_0)$ ,  $Q_1 = [\alpha_1, \alpha_2)$ ,  $\dots$ ,  $Q_r = [\alpha_r, \infty)$ , such that, for each  $j \in \{0, 1, \dots, r\}$ ,  $S_j$  is the unique maximal minimizer of  $h_\alpha(S)$ ,  $\forall \alpha \in Q_j$ . A minimizer of Eq. (12)  $\mathbf{t}^* = (t_1^*, t_2^*, \dots, t_{|V|}^*)$  is then determined:  $t_i^* = \frac{f(S_{j+1}) - f(S_j)}{\mathbf{1}(S_{j+1} \setminus S_j)}$ ,  $\forall i \in (S_{j+1} \setminus S_j)$ ,  $j = (1, \dots, r)$ . We introduce two lemmas [20] to see that  $l$  is a non-decreasing submodular function.

**Lemma 1 (Lemma).** *For any  $\eta \in \mathbb{R}$  and submodular function  $h$ ,  $\mathbf{t}^*$  is an optimal solution to  $\min_{\mathbf{t} \in \mathcal{B}(l)} \|\mathbf{t}\|_2^2$  if and only if  $\mathbf{t}^* - \eta \mathbf{1}$  is an optimal solution to  $\min_{\mathbf{t} \in \mathcal{B}(l) + \eta \mathbf{1}} \|\mathbf{t}\|_2^2$ .*

**Lemma 2 (Lemma).** *Set  $\eta = \max_{i=1, \dots, |V|} \{0, l(V \setminus \{i\}) - l(V)\}$ , and then  $l + \eta \mathbf{1}$  is a non-decreasing submodular function.*

With Lemma 2, we solve

$$\min_{S \subset V} f(S) - \hat{z}(S) + (\eta - \alpha) \mathbf{1}(S), \quad (15)$$

and apply Lemma 1 to obtain a solution to the original problem. With fixed  $\alpha$ , we can efficiently attain the optimal of Eq. (15) because this is a minimum cut problem.

**Proposition 1.** *The problem in Eq. (15) is equivalent to a minimum s/t-cut problem.*

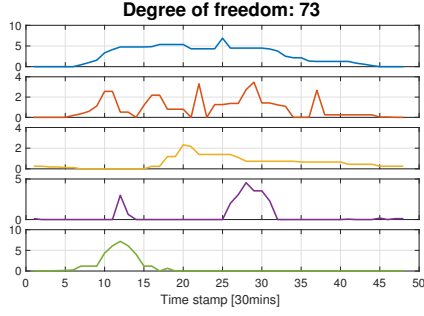


Fig. 2: Time factors of Proposed 1 on the traffic flow data

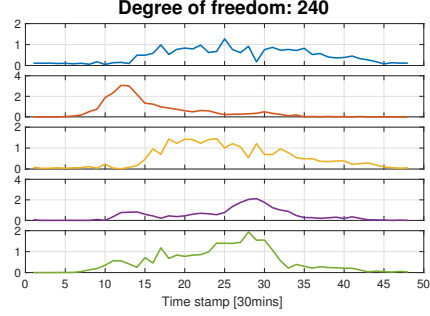


Fig. 3: Time factors of Multi+Lap on the traffic flow data

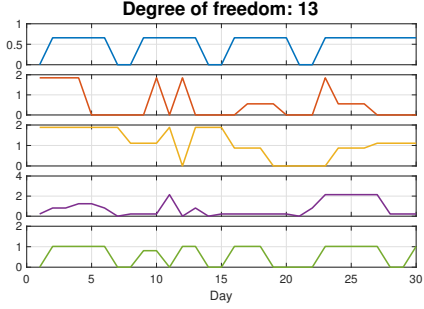


Fig. 4: Day factors of Proposed 1 on the traffic flow data

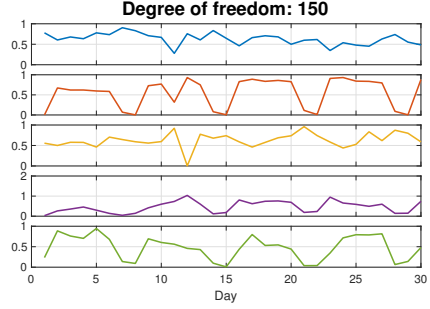


Fig. 5: Day factors of Multi+Lap on the traffic flow data

*Proof.* Each component in  $f$  is graph-representable. The graph is obtained due to the additive property of the graph-representative submodular functions, where the groups of parameters are represented with hyper nodes  $u_1^k, u_0^k$  that corresponds to each group, and the capacities of the edges between hyper and ordinal nodes  $v_i \in V$ .

The attained graph includes both of the GFL and HOFL graphs as spacial cases. As a consequence, we can attain a sequence of solutions for all  $\alpha$  of the parametric  $s/t$  minimum-cut problem (15) using an efficient parametric-flow algorithm, such as [12], that runs in  $O(|V'| |E'| \log(|V'|^2 / |E'|))$  as the worst case and  $|V'|$  and  $|E'|$  are the number of nodes and edges of the graph.

## References

1. Muhammad Intizar Ali, Feng Gao, and Alessandra Mileo. Citybench: A configurable benchmark to evaluate rsp engines using smart city datasets. In *ISWC*, pages 374–389, 2015.
2. Francis R Bach. Structured sparsity-inducing norms through submodular functions. In *Proc. of NIPS*, pages 118–126, 2010.
3. Alvaro Barbero and Suvrit Sra. Fast newton-type methods for total variation regularization. In *Proc. of ICML*, pages 313–320, 2011.

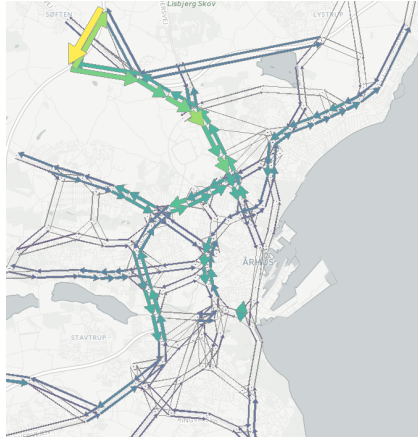


Fig. 6: A spatial pattern of the blue factor of Proposed 1 on the traffic flow data

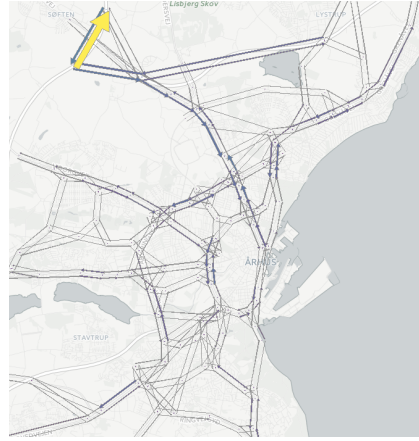


Fig. 7: A spatial pattern of the blue factor of Multi+Lap on the traffic flow data

4. Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
5. Yuri Boykov and Vladimir Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *IEEE transactions on pattern analysis and machine intelligence*, 26(9):1124–1137, 2004.
6. Deng Cai, Xiaofei He, Jiawei Han, and Thomas S Huang. Graph regularized nonnegative matrix factorization for data representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1548–1560, 2011.
7. Antonin Chambolle and Jérôme Darbon. On total variation minimization and surface evolution using parametric maximum flows. *International journal of computer vision*, 84(3):288, 2009.
8. Andrzej Cichocki, Rafal Zdunek, Anh Huy Phan, and Shun-ichi Amari. *Nonnegative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*. John Wiley & Sons, 2009.
9. Inderjit S Dhillon and Suvrit Sra. Generalized nonnegative matrix approximations with bregman divergences. In *Proc. of NIPS*, volume 18, 2005.
10. Zipei Fan, Xuan Song, and Ryosuke Shibasaki. Cityspectrum: a non-negative tensor factorization approach. In *Proc. of UbiComp*, pages 213–223, 2014.
11. Satoru Fujishige. *Submodular functions and optimization*, volume 58. Elsevier, 2005.
12. Giorgio Gallo, Michael D Grigoriadis, and Robert E Tarjan. A fast parametric maximum flow algorithm and applications. *SIAM Journal on Computing*, 18(1):30–55, 1989.
13. Nicolas Gillis. The why and how of nonnegative matrix factorization. In *Regularization, Optimization, Kernels, and Support Vector Machines*, pages 257–291. Chapman and Hall/CRC, 2014.
14. Yufei Han and Fabien Moutarde. Analysis of large-scale traffic dynamics in an urban transportation network using non-negative tensor factorization. *International Journal of Intelligent Transportation Systems Research*, 14(1):36–49, 2016.
15. Kejun Huang, Nicholas D Sidiropoulos, and Athanasios P Liavas. A flexible and efficient algorithmic framework for constrained matrix and tensor factorization. *IEEE Transactions on Signal Processing*, 64(19):5052–5065, 2016.

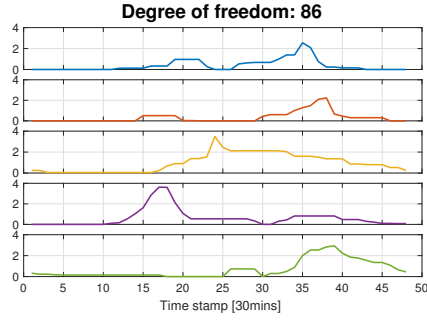


Fig. 8: Time factors of Proposed 1 on the bike-sharing data of Washington D.C.

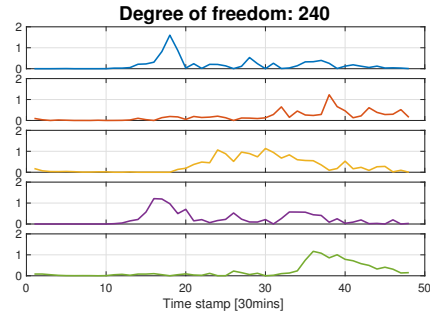


Fig. 9: Time factors of Multi+Lap on the bike-sharing data of Washington D.C.

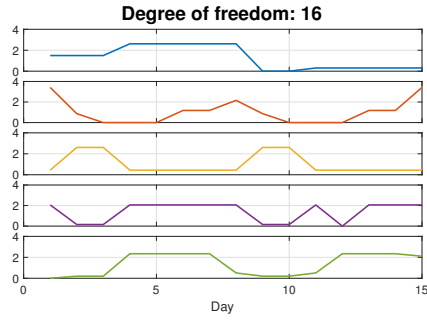


Fig. 10: Day factors of Proposed 1 on the bike-sharing data of Washington D.C.

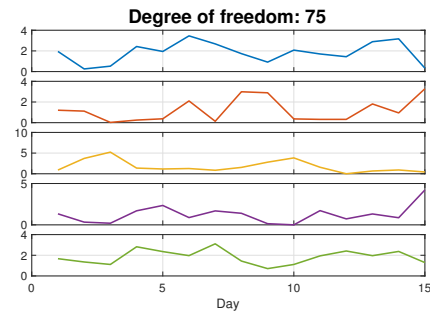


Fig. 11: Day factors of Multi+Lap on the bike-sharing data of Washington D.C.

16. Tatsuaki Kimura, Keisuke Ishibashi, Tatsuya Mori, Hiroshi Sawada, Tsuyoshi Toyono, Ken Nishimatsu, Akio Watanabe, Akihiro Shimoda, and Kohei Shiomoto. Spatio-temporal factorization of log data for understanding network events. In *Proc. of INFOCOM*, pages 610–618, 2014.
17. Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
18. Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
19. László Lovász. Submodular functions and convexity. In *Mathematical Programming The State of the Art*, pages 235–257. Springer, 1983.
20. Kiyohito Nagano, Yoshinobu Kawahara, and Kazuyuki Aihara. Size-constrained submodular minimization through minimum norm base. In *Proc. of ICML*, pages 977–984, 2011.
21. Nikhil Rao, Hsiang-Fu Yu, Pradeep K Ravikumar, and Inderjit S Dhillon. Collaborative filtering with graph information: Consistency and scalable methods. In *Proc. of NIPS*, pages 2107–2115, 2015.
22. Alexander J Smola and Risi Kondor. Kernels and regularization on graphs. In *Learning theory and kernel machines*, pages 144–158. Springer, 2003.
23. Lijun Sun and Kay W Axhausen. Understanding urban mobility patterns with a probabilistic tensor factorization framework. *Transportation Research Part B: Methodological*, 91:511–524, 2016.

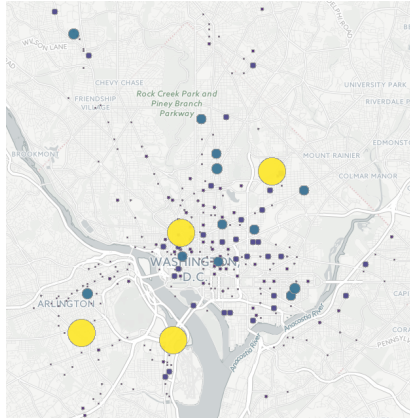


Fig. 12: A spatial pattern of the yellow factor of Proposed 1 on the bike-sharing data of Washington D.C.

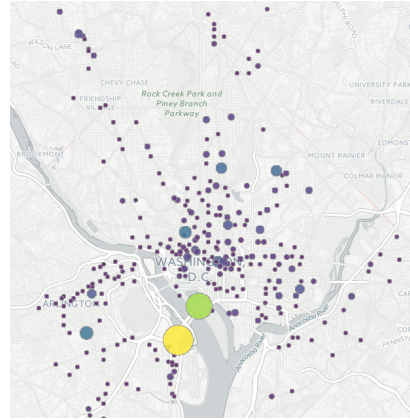


Fig. 13: A spatial pattern of the yellow factor of Multi+Lap on the bike-sharing data of Washington D.C.

24. Koh Takeuchi, Yoshinobu Kawahara, and Tomoharu Iwata. Higher order fused regularization for supervised learning with grouped parameters. In *Proc. of ECMLPKDD*, pages 577–593, 2015.
25. Koh Takeuchi, Ryota Tomioka, Katsuhiko Ishiguro, Akisato Kimura, and Hiroshi Sawada. Non-negative multiple tensor factorization. In *Proc. of ICDM*, pages 1199–1204, 2013.
26. Rashish Tandon and Suvrit Sra. Sparse nonnegative matrix approximation: new formulations and algorithms. *Rapport technique*, 193:38–42, 2010.
27. Robert Tibshirani, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1):91–108, 2005.
28. Yu-Xiang Wang, James Sharpnack, Alex Smola, and Ryan J Tibshirani. Trend filtering on graphs. *Journal of Machine Learning Research*, 17(105):1–41, 2016.
29. Daniela M Witten, Robert Tibshirani, and Trevor Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, page kxp008, 2009.
30. B. Xin, Y. Kawahara, Y. Wang, and W. Gao. Efficient generalized fused lasso with its application to the diagnosis of alzheimer’s disease. In *Proc. of AAAI*, pages 2163–2169, 2014.
31. Yangyang Xu and Wotao Yin. A block coordinate descent method for regularized multi-convex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences*, 6(3):1758–1789, 2013.
32. Dingqi Yang, Daqing Zhang, Vincent W Zheng, and Zhiyong Yu. Modeling user activity preference by leveraging user spatial temporal characteristics in lbsns. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 45(1):129–142, 2015.

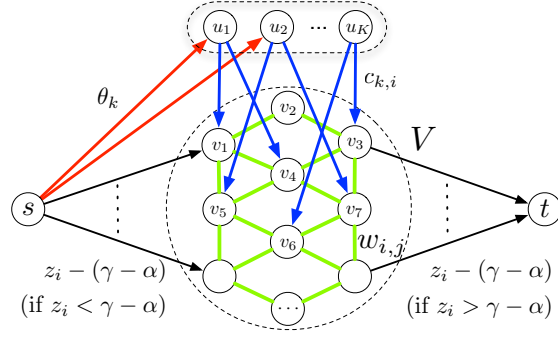


Fig. 14: Minimum  $s/t$ -cut problem of Problem (15). Given graph  $G = (V', E')$  for our proposed method, capacities of edges  $c(v'_i, v'_j)$  are defined as:  $c(s, u_k) = \theta_k$ ,  $c(v_i, v_j) = w_{i,j}$ ,  $c(u_k, v_i) = c_{k,i}$ ,  $c(s, v_i) = z_i - (\gamma - \alpha)$  if  $z_i > \gamma - \alpha$ , and  $c(v_i, t) = (\gamma - \alpha) - z_i$  if  $z_i < \gamma - \alpha$ . Nodes  $u_k$   $k = (1, \dots, K)$  are hyper nodes that correspond to the groups  $g_k$ . And  $s$ ,  $t$ , and  $v_i$  are source, sink, and parameters nodes, respectively.